

Parallel Discrete Event Simulation Course #9

David Jefferson
Lawrence Livermore National Laboratory
2014

This work was performed under the auspices of the U.S. Department
of Energy by Lawrence Livermore National Laboratory under Contract
DE-AC52-07NA27344. Lawrence Livermore National Security, LLC

Release Number: LLNL-PRES-652636

Reprise

Classic Time Warp Algorithm: Local Synchronization

Parallel Discrete Event Simulation -- (c) David Jefferson, 2014

3

We call this the “classic” Time Warp algorithm because it most closely resembles the very first implementation. Current implementations differ in many ways, though not fundamentally. This presentation will also emphasize clarity and abstraction and symmetry in the algorithm, so that it is easiest to explain, even if actual implementations differ in many details.

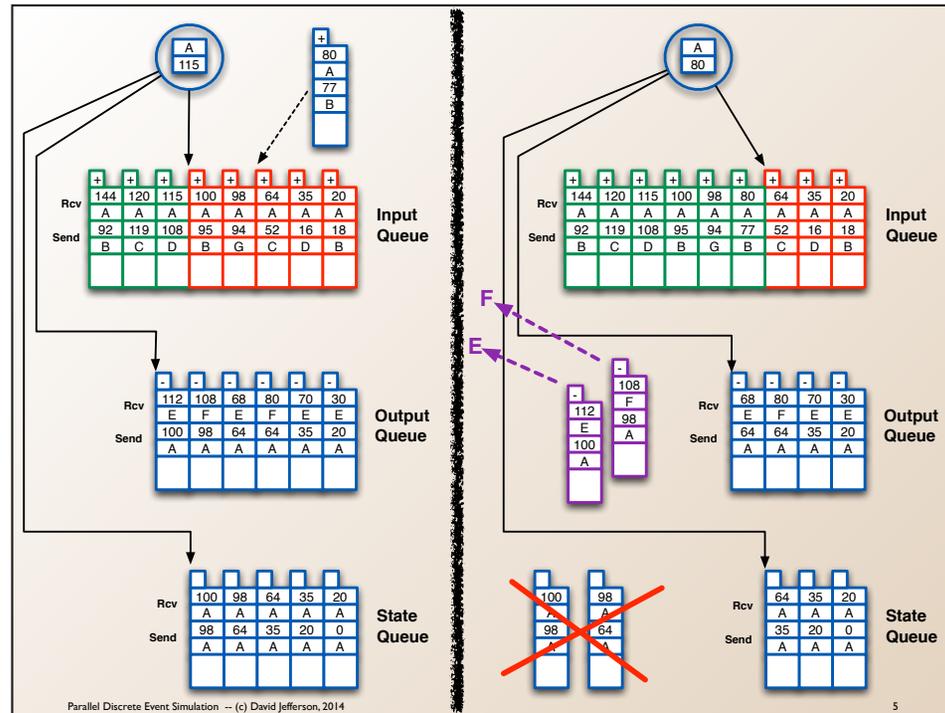
Asynchronous, distributed rollback? Are you serious???

- **Must be able to restore any previous state (between events)**
- **Must be able to cancel the effects of all “incorrect” event messages that should not have been sent**
 - even though they may be in flight
 - or may have been processed and caused other incorrect messages to be sent
 - to any depth
 - including cycles back to the originator of the first incorrect message!
- **Must do it all *asynchronously*, with *many concurrently interacting rollbacks in progress*, and *without barriers***
- **Must deal with the consequences of executing events starting in “incorrect” states**
 - runtime errors
 - infinite loops
- **Must guarantee global progress (if sequential execution progresses)**
- **Must deal with truly irreversible operations**
 - I/O, or freeing storage, or launching a missile
- **Must be able to operate in finite storage**
- **Must achieve good parallelism, and scalability**

Parallel Discrete Event Simulation -- (c) David Jefferson, 2014

4

This shows the list of challenges we have to overcome for the Time Warp algorithm, or any optimistic PDES algorithm, to be practical. Most people with a background in asynchronous distributed computation who have not seen optimistic PDES algorithms are inclined to believe that doing this is either literally impossible, or at least hopelessly complex and slow.



This diagram represents the before and after of the arrival of a straggler message.

Object A is at simTime 115 when an event message arrives in the past with a timestamp of 80. Two events at time 98 and 100 should not have been executed. The event at time 115 has been partially executed--we are still in the middle of executing it--and it may have sent some but not all of the event messages it would have. In any case, whatever it did is likely to be wrong.

We have to roll back to the state at time 80, i.e. to the state saved after execution of the event at time 64, the last correctly-executed event before time 80. The rollback consists of the following actions. Strictly speaking the kind of rollback we are describing here is called "aggressive cancellation", and is in some ways the most "optimistic" of the optimistic algorithms. An alternative is called "lazy cancellation", and another is called ??

The rollback consists of the following steps.

- 1) Insert the straggler message into the input queue where it belongs in the sort. It may be an antimessage and annihilate with another message already in the input queue. That make no difference in the algorithm at all--antimessages are treated identically.
- 2) Interrupt the event in progress (115). Restore the state saved after event 64 as the current state of the object. Delete the two subsequently saved states created by events 98 and 100, since they are (probably) incorrect (and one of them has been partially modified by event 115). Note that in a rollback variation called "lazy re-evaluation" these states would not actually be deleted at this time--and hence it really is possible to have "future" states in the queue. And in another variation called "sparse state saving" where we don't save states between every two event, but do it less often than that, then we might have to roll back farther than time 80 restore to an even earlier state.)
- 3) From the output queue, find the antimessages to the messages sent incorrectly after time 80, dequeue them, and deliver them to their receivers--the same objects that the original (incorrect) positive messages went. Note that this includes any messages sent by the event that was in progress (and was interrupted), in this case event 115. **Surprisingly (!) that is all that is required to exactly undo the effects of those positive messages, whether they have been delivered yet or not, or have been processed or not, or have caused generation of a tree of further, probably incorrect, distributed computation.** The fact that this antimessage mechanism works in all cases, and allows the simulation globally to make progress asynchronously, independent of the speed of execution of the objects or the latency of message delivery, and regardless of the possibility of many interacting rollbacks in progress simultaneously, is a key observation at the foundation of most optimistic methods. (However, some less aggressive variations, e.g. *risk free* algorithms (in Paul Reynolds' taxonomy) do not transmit event messages until they can be committed, and thus have no need for antimessages. This comment is a forward reference, and I don't know whether I will get back to it in the course.)

Global Virtual Time (GVT) and Commitment: Global Synchronization

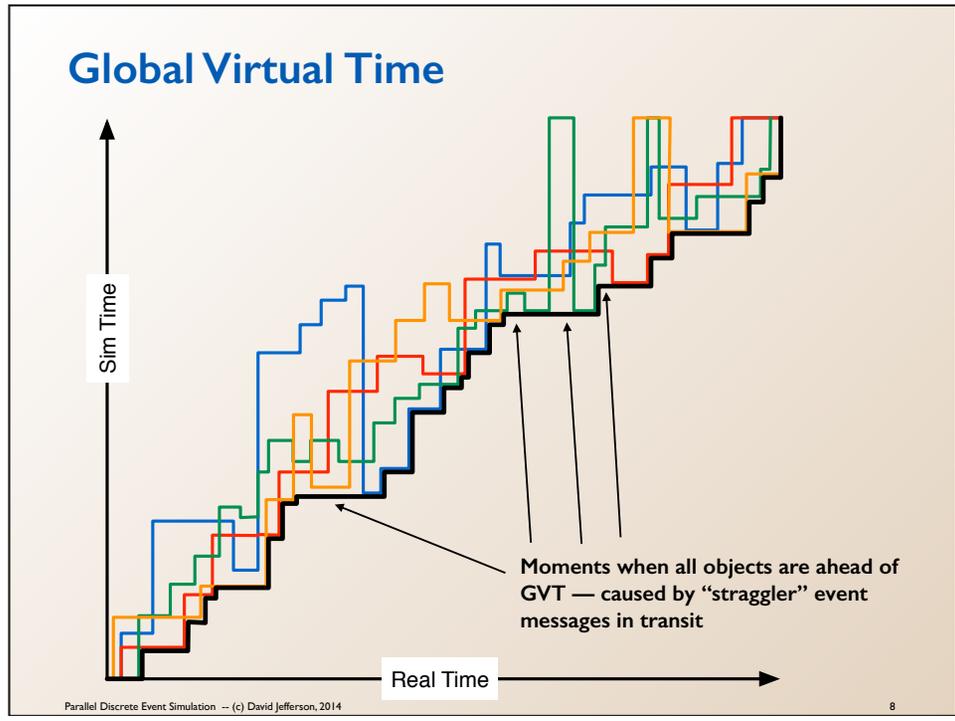
Local Virtual Time (LVT) and Global Virtual Time (GVT)

- **Virtual times related to, but not the same as, simulation time.**
 - In the literature we have sometimes defined *simulation time* to be just the high order bits of *virtual time*, but virtual time is a wide “address in time”.
 - We also use the term “virtual time” because of a strong analogy to “virtual memory (to be developed later)”
- **The Local Virtual Time (LVT) of an object measures how far that object has progressed simulation time, i.e. what its simulation clock reads.**
 - If an object is blocked because it has (temporarily) executed all of the events in its input queue, then we define its LVT = ∞
- **Global Virtual Time (GVT) measures how far the entire simulation has progressed globally, and is (roughly) the minimum of all of the LVTs.**

Parallel Discrete Event Simulation -- (c) David Jefferson, 2014

7

LVTs go forward and backward in time, but more often forward. GVT *never* goes backward, and is *always* (at any instant of wall clock time) a lower bound for all LVTs at that instant.



There are a few places in this diagram where GVT is strictly lower than the minimum of all LVTs. That will happen at times when a message is in transit that happens to carry a lower timestamp (receive time) than the LVT of any object. When it is delivered, it will cause the receiving object to rollback to the message's received time.

Properties of Global Virtual Time

- **Events at virtual times lower than GVT can never be rolled back.**
- **GVT never decreases.**
 - In a well-posed simulation GVT inevitably *increases*.
- **GVT == ∞ is criterion for “normal” termination**

Parallel Discrete Event Simulation – (c) David Jefferson, 2014

9

By definition, GVT never decreases, and in fact it must increase unless the simulation has a bug in it like an infinite loop that would also affect the sequential execution of the same model.

GVT == ∞ if and only if all objects are at time infinity and no messages are in transit. In that case all objects have run out of events to process, and the entire global simulation terminates normally. Termination detection is the same as detecting the GVT == ∞ .

Definition of *instantaneous* Global Virtual Time

$$GVT = \min (LVT(p), RT(q), ST(r))$$

objects: p
forward messages in transit from p: q
reverse messages in transit from p: r

Parallel Discrete Event Simulation -- (c) David Jefferson, 2014

10

LVT(p) = Local Virtual Time, i.e. the simulation clock value of Object p

RT(q) = Receive Time of the (positive or negative) event message q that is in transit

ST(r) = Send Time of (positive or negative) event message r that is in transit in the reverse direction from receiver's output queue to sender's input queue (for storage management/ flow control)

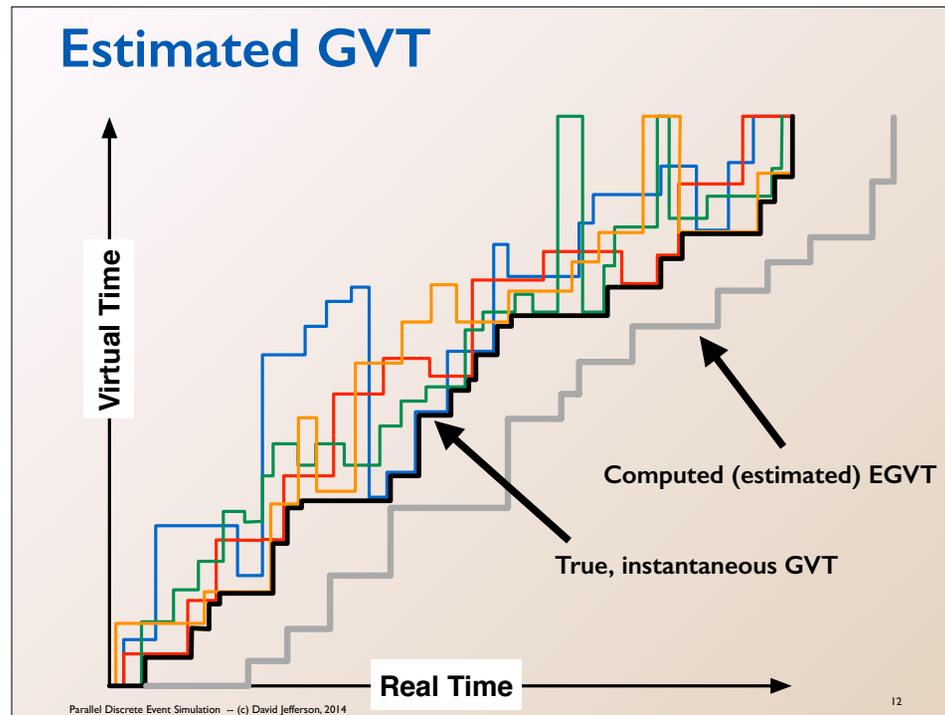
This is an *instantaneous* definition. It could be only calculated exactly if we stopped the simulation globally and waited for delivery of all messages. However, in practice, we calculate an *estimate* of it asynchronously, without a barrier, while objects continue executing and messages continue to be transmitted.

At least half a dozen algorithms for estimating GVT and broadcasting the result without any barrier synchronization have been published. All take time $O(\log n)$ where n is the number of processes. They take advantage of the fact that a message may be in transit if it has been sent but not acknowledged yet (in a low-level reliable transmission protocol). Just because a message has not yet been acknowledged, that does not mean it has not actually been delivered yet--it may have been, but the ack has not yet arrived. In that case the message will be included in the "min" operation of both sender and receiver. But that does not change the estimated GVT value.

The estimate is guaranteed to be low, which is the direction you want it to be. An estimate that is too high would cause Time Warp to commit events that are not yet safe from rollback--that would be a disaster. But an estimate that is too low just delays the commitment of some events that are in fact safe to commit.

Estimated GVT

- **Estimated periodically and also when memory is exhausted**
- **Value is broadcast to all objects**
- **Objects then locally perform commitment operations and storage recovery**
- **Same quantity (GVT) is used for virtually all conservative algorithms, but used in different ways**
 - **safety for conservative algorithms**
 - **commitment and storage management for optimistic algorithms**



True, instantaneous GVT cannot be calculated without a barrier pausing the simulation which, besides yielding very poor performance as the simulation parallelism declines to zero, is also difficult to schedule if out of memory, and also unnecessary. We want to calculate GVT asynchronously, and only when necessary, without pausing event execution.

Instead of calculating true GVT, we calculate an *estimate* of GVT, EGVT, both periodically. The estimate needs to satisfy two key properties:

- 1) It is always less than or equal to true instantaneous GVT, i.e. it is a lower bound on true GVT.
- 2) It “tracks” true GVT, never very far behind true GVT.

Regarding the second condition, the definition of “tracking” is not obvious. We do not mean that it stays within a constant difference or constant ratio of numeric value of GVT. What we mean is that at the time a new value of EGVT is computed, EGVT is a value that GVT exceeded no longer ago than a constant amount of real time earlier. Thus EGVT is an “out of date” value of GVT, a value that GVT exceeded no more than a short time earlier. How short a time? The time it takes to compute EGVT! In other words, EGVT is greater than or equal to the value that GVT has at the start of the computation of EGVT, but less than or equal to the value GVT has when EGVT calculation is complete.

Commitment

- Commitment means *giving up the option to roll back*.
- Since we know that no object will ever have to roll back to any simulation time $< \text{GVT}$, then we can *commit* all events and release all resources related to simtimes that are $< \text{GVT}$
- We calculate EGVT periodically, and thus commit periodically

EGVT calculation

- We can use an *estimate* of GVT, a recent true value of GVT, rather than the current instantaneous value
 - Such an estimate must never be high
 - But it should not be too far out of date either
- Calculated EGVT periodically, or sooner if memory is exhausted on some node
- EGVT can be calculated asynchronously, while simulation continues, without barriers
- EGVT is broadcast to all objects
- Objects then locally perform commitment operations and storage recovery
- Same quantity is used for virtually all conservative algorithms, but used in different ways
 - *safety* for conservative algorithms
 - *commitment* for optimistic algorithms

GVT estimation algorithms (EGVT)

- **Many algorithmic variations published in the literature**
 - Synchronous and asynchronous
 - One-sided and two-sided
 - Require end-to-end message acks or not
- **The most useful ones have these properties**
 - Execute concurrently with event execution
 - Take $(\log n)$ -time (assuming constant message latency)
- **For a survey see**
 - Fujimoto, [Parallel and Distributed Simulation Systems](#), Wiley, 2000, Section 4.4*
 - or look up “simulation GVT calculation” in Google

Uses of EGVT: Commitment actions

- **“Fossil collection”**
 - free the memory for input messages, output messages, and states no longer needed to support rollback
- **Termination detection**
 - check whether $GVT == \infty$
- **I/O commitment**
 - Irreversible output operations at time before GVT that were postponed can now be committed (in increasing virtual time order)
 - Input operations done before GVT for which the option to “un-input” was preserved can now be committed, and any buffers freed.
- **Runtime error handling**
 - All runtime errors must be trapped by the simulator
 - Saved states should be marked as to whether or not a runtime error occurred during their production (and if so, what the error was)
 - If any saved state that is marked in error is committed, then the whole simulation was in error, and must be terminated.
- **Event message transmission in “risk free” Time Warp variation**
 - Outgoing event messages can be delayed until commit time before being transmitted
 - Then there is no “risk” that they will need to be cancelled, so no output queue is necessary and no negative messages need to be created

Parallel Discrete Event Simulation -- (c) David Jefferson, 2014

16

“Fossil collection” is supposed to remind you of “garbage collection”, i.e. the recycling of storage that can never be accessed again. Events that are for time lower than EGVT can never be rolled back, so it is safe to release the memory that has been used to hold the input messages, output messages, and states associated with those states.

We already mentioned that $EGVT == \infty$ equivalent to global normal termination. Abnormal termination, however, generally occurs at some finite simulation time. Sometimes when the simulator is told to cut off execution artificially at some particular simulation time t_{max} , that is implemented by cutting it off whenever GVT is first calculated to be $\geq t_{max}$.

Input is generally handled by buffering input data so that it can be “un-input” in the case of a rollback. The data buffers are released once GVT has increased past the time of the event that did the inputting, to where the data will never have to be un-input.

Output is handled by buffering the data until such time as the event requesting the output is committed. Then the data can be physically written out, and will never have to be un-output.

If an optimistic simulation needs to both read and write the same file or database, then that file or data base has to be treated as one or more full-fledged Objects in the simulation, and it need to have a simulation clock (virtual time) and it must be able to roll back.

Asynchronous, distributed rollback? Are you serious???

- Must be able to restore any previous state (between events)
- Must be able to cancel the effects of all “incorrect” event messages that should not have been sent
 - even though they may be in flight
 - or may have been processed and caused other incorrect messages to be sent
 - to any depth
 - including cycles back to the originator of the first incorrect message!
- Must do it all *asynchronously, with many concurrently interacting rollbacks in progress, and without barriers*
- Must deal with the consequences of executing events starting in “incorrect” states
 - runtime errors
 - infinite loops
- Must guarantee global progress (if sequential model progresses)
- Must deal with truly irreversible operations
 - I/O, or freeing storage, or launching a missile
- Must be able to operate in finite storage
- Must achieve good parallelism, and scalability

Parallel Discrete Event Simulation – (c) David Jefferson, 2014

17

This shows the list of challenges we have to overcome for the Time Warp algorithm, or any optimistic PDES algorithm, to be practical.

In green I have listed those that we have covered so far

End Reprise

A collection of remaining issues

Other issues

- Runtime Errors
- Infinite Loops
- GVT estimation algorithms
- Scheduling
- Shared Memory implementations
- Multicore implementations
- Flow Control
- Storage Management
- Throttling
- TW as an operating system
- Variations on rollback
- Variations on message cancellation
- Mixed conservative and optimistic synchronization
- Checkpoint / Restart
- Reverse computation
- Measured performance at extreme scale
- Fault recovery
- Load balancing
- Symmetry
- Applications beyond PDES

Runtime errors

- During a “correct” event:
 - Runtime error caught
 - Will behave the same way as it would in the sequential algorithm
 - Runtime error uncaught (e.g. errant pointer)
 - Both optimistic and sequential algorithms may misbehave, but possibly differently.
- During an “incorrect” event, executed from an “incorrect” state
 - Runtime error caught
 - TW should stop executing the event, mark the object as being “Erroneous”, and “output” a tentative error message
 - If the event is rolled back, remove the “Erroneous” mark, and proceed normally
 - If an Erroneous event is committed, then the error message is output and the simulation should be terminated globally
 - Runtime error uncaught (e.g. errant pointer)
 - Program will misbehave, *and that misbehavior may not be corrected by subsequent rollback!*
 - Hence, we really need a *proof* that an event method cannot cause an uncaught error from *any* state before we can safely execute it optimistically.

Infinite loops within events

- During a “correct” event:
 - The event behaves the same way it would in the sequential algorithm
 - The infinite loop is then “correct”
 - GVT never increases
 - Tough!
- During an “incorrect” event, executed from an “incorrect” state
 - If events are interrupted when stragglers arrive that would cause a rollback
 - Eventually a straggler will arrive (guaranteed) which will cancel this event or require a rollback to an earlier state.
 - Execution will proceed properly and duplicate behavior of the sequential algorithm
 - If events are executed uninterruptably
 - The infinite loop is never interrupted
 - GVT never increases
 - **This is *not* what the sequential algorithm does!**
 - Hence, we really need a *proof* that an event method cannot cause an infinite loop from *any* state before we can safely execute it optimistically.

Scheduling

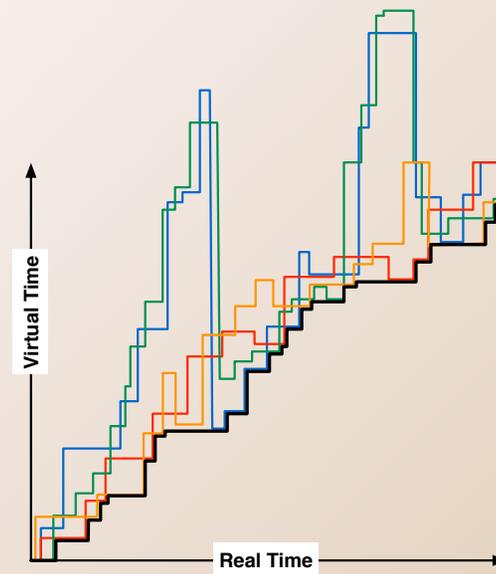
- **Assume**
 - k cores per compute node
 - multiple (sequential) processes or tasks per compute node
 - multiple objects or LPs per process or task
 - Each process (task) can be viewed as a “fat object”
 - LVT of a process (task) is the minimum LVT of any object in it
 - A process is only blocked ($LVT = \infty$) if all objects in it are blocked
- **Natural scheduling discipline for processes in shared memory and objects within processes, is LVTF (lowest virtual time first)**
 - System functions (EGVT calculation, storage reclamation) prioritized over event execution
 - On each compute node always run the k (non-blocked) processes with lowest LVTs
 - Within each process, run the one object with lowest LVT
 - Preemptive scheduling within events preferred, but not always implementable

Shared memory and / or multicore architectures offer significant optimizations

- **If sender and receiver of message are in the same memory**
 - There can be a single buffer holding the message, with + and - pointers to them
 - Determination of whether two messages are antimessages is just a pointer comparison
 - Messages “delivered” FIFO, i.e. order preserved in real time
- **Each process can be treated as a single “fat” LVT**
 - Speeds EGVt calculation
- **Shared event message queues**
 - Shared high-level event queue among processes
 - Shared lower level queues within processes

Throttling excess optimism

- Conservative algorithms often make slower progress than they should -- they are too conservative.
- Optimistic algorithms often progress (locally) faster than they should, and then have to roll back a lot of work -- they can be too optimistic.
- The costs of excess optimism
- Throttling is any mechanism that reduces excess optimism
- Throttling algorithms in the literature
 - Moving Time Windows (w)
 - Breathing Time Warp



Parallel Discrete Event Simulation -- (c) David Jefferson, 2014

25

What is the cost of excess optimism?

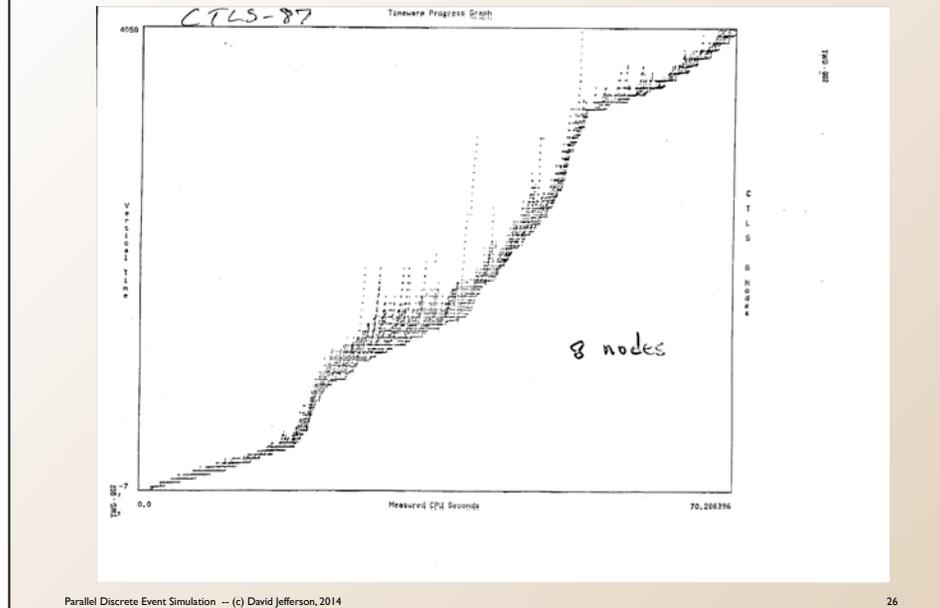
- 1) Extra storage impact
- 2) The cost of local rollback, which may be mostly storage reclamation, but in some algorithms (like reverse computation) may require compute time proportional to that used in forward execution.
- 3) Overhead costs in other processes for dealing with messages that should not have been sent (unless conservative message transmission was used).

Moving time windows: This throttling protocol prevents runaway execution into the future by setting a simulation time window size w . The protocol requires an Object to block once it reaches time $(GVT + w)$, i.e. there is a global barrier at $(GVT+w)$. Once all processes have reached that time, then GVT is re-estimated, and a new $(GVT+w)$ window is enforced.

Window size w can, of course, change dynamically and adaptively. But what value should w have? And when should it be increased or decreased, and by how much? Today this is handled empirically. Without a good theory for adaptively setting w , I find it unsatisfactory.

Variation: The window size can be in number of events instead of units of simulation time. That has the virtue that it does not use the arithmetic properties of simtime, but only the total ordering properties. But otherwise the same issues arise as when the window size is in units of simulation time: how big should the window be, and how should it be increased and decreased.

Tactical military simulation II on 8 nodes



This is a trace of a different (and more complex) tactical force-on-force military simulation. It was run on only 8 nodes, and while the speed of progress through virtual time varies considerably, as indicated by the slope of GVT, the progress was relatively tight, with only a few occasions where an object ran forward ahead of GVT and eventually rolled back. *But compare this to the next slide.*

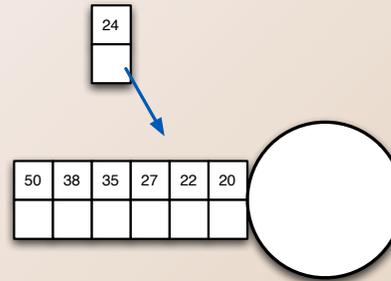
This data was collected at JPL circa 1990 on the 64-node Caltech Hypercube (Motorola 68020, 4 MB per node) running the Time Warp Operating System. Sponsored by TRADOC, U.S. Army.

Throttling algorithms

- **MTW: Moving Time Windows**
 - Establish a “window” of width w units of virtual time
 - Do not allow any object to execute beyond $(EGVT + w)$
 - The window bound $(EGVT + w)$ is advanced every time $EGVT$ is re-calculated
 - The window width w can change dynamically if desired
 - **Problems:**
 - Why should w be global?
 - What is the basis for deciding on w ?
 - Blocks a process even if it is making correct progress
- **BTW: Breathing Time Warp**
 - Similar to MTW, except that the window width is measured in a certain number of events, instead of an amount of virtual time

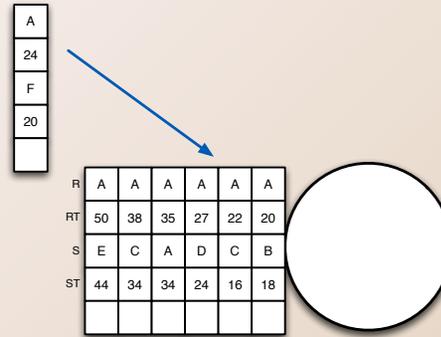
Message flow control in PDES

- **Flow control** is the memory management problem that arises when producers send data faster than consumers process it.
- For graph-oriented, conservative PDES protocols with FIFO order preservation on each channel then *windowing protocols* on a per-channel basis work fine.
- But windowing does not work for optimistic PDES:
 - No channelization, so windowing does not apply
 - Messages not transmitted or processed in FIFO order!
 - Messages not immediately deleted even after they are processed!
 - When memory is full, a message with a low time stamp may still arrive and there may be no room for it!



Optimistic flow control

- Suppose a message from F arrives at A when A's memory is full? What do we do?
- We cannot just block the sender, F — it is already too late!
- We cannot just drop the message, because it has to be processed before some others we have enqueued.
- EGVT update may release some space, but it is not guaranteed to.
- Answer: *Send back the highest RT message from A to E, in this case causing E to roll back to 44, before E sent it.*
- This is a message sent *backward in space* and *backward in virtual time!*
- In this protocol a message from F to A causes E to roll back!



Cancelback protocol: Generalization to all dynamic storage

- Whenever memory has been fossil-collected, storage is full, and a new item (state, input msg, output msg) has to be allocated:

Find *any* item in memory (including the new item) with $ST > EGV$ (preferably with maximal ST):

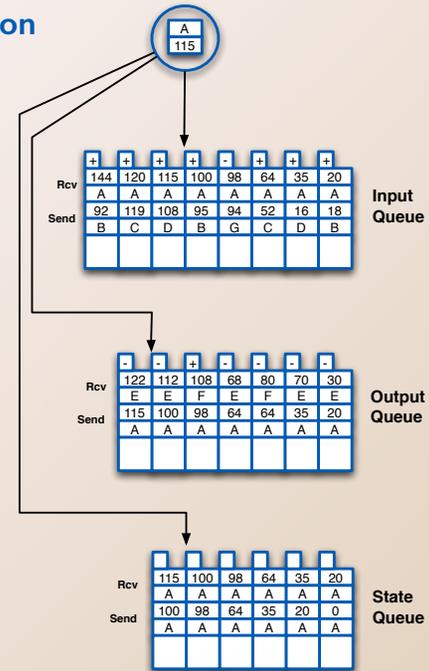
if none: **fail: out of memory**
 if item is a state: delete it
 if item is an output message: send it forward
 if item is an input message: send it backward
 Roll back "victim" object if necessary

- Note: Item in memory can be found in a queue of another *completely unrelated object* in the same memory pool.

- Provably space optimal (but not time optimal)!

- See

Jefferson, David, "Virtual Time II: The Cancelback Protocol for Storage Management in Time Warp", Proc. ACM Symp. on the Princ. of Distributed Computing (PODC), Quebec City, Quebec, August 1990



Time Warp as an operating system

- **temporal coordinate system:**
 - nothing comparable in OS world
- **process name space:**
 - flat space, like MPI, unlike other OSs
- **process (LP) scheduling:**
 - LVTF - lowest virtual time first
- **synchronization:**
 - virtual time, rollback
- **determinism:**
 - guaranteed, unlike with most OSs
- **interprocess communication:**
 - timestamped event messages to named objects
- **queueing:**
 - virtual time order, with antimessage annihilation
- **flow control:**
 - message sendback
- **memory management:**
 - fossil collection of past, throttling or cancel back of future
- **delayed commitment**
 - tied to EGVF
 - nothing comparable in OS world
- **memory alloc / dealloc:**
 - freeing tied to commitment
- **error handling:**
 - tied to commitment
- **I/O:**
 - tied to commitment
- **job termination:**
 - tied to commitment

Parallel Discrete Event Simulation -- (c) David Jefferson, 2014

32

One way to view Time Warp is as a special purpose parallel operating system. It can be viewed as having all of the architecture of a parallel OS, but different algorithms in place of most of the OS components.

- 1) *Temporal coordinate system:* TW offers a temporal coordinate system (address space), namely virtual time, which has no analog in other OSs.
- 2) *Process name space:* TW has a flat, global process name space not offered by most parallel OSs, although it is offered by MPI.
- 3) *Process scheduling:* TW generally uses LVTF (lowest virtual time first) scheduling among logical processes. The only comparable scheduling algorithm used in conventional OSs is a strict priority-driven scheduling. LVTF can be replaced by any scheduling mechanism that eventually executes the lowest
- 4) *Synchronization:* TW uses virtual time and rollback (with some process blocking) as its fundamental synchronization mechanism. Other OSs use various kinds of process blocking mechanisms in the context of locks, semaphores, messages waiting, etc.
- 5) *determinism:* TW guarantees deterministic execution; other OSs do not
- 6) *interprocess communication:* TW uses event messages signs and both sender and receiver event stamps; other OSs have a variety of mechanisms for interprocess communication, including shared memory, pipes, TCP/IP, signals, MPI, etc.
- 7) *queueing:* TW queueing is a priority queue with virtual Time as priority, but with the added feature of antimessage cancellation. Most message queues in other OSs are FIFO.
- 8) *flow control:* Flow control in TW can use message sendback, whereas that is not possible in other OSs because the sending processes cannot be rolled back.
- 9) *memory management:* process memory is managed by a combination of fossil collection (reclamation of space with time stamps in the past) and cancel back (reclamation of space in the future). Nothing like it exists for conventional OSs.
- 10) *delayed commitment:* in TW commitment of irreversible operations is delayed until commitment time. In most OSs, since there is no rollback, and commitment time is immediate. (Notable exceptions are DB management systems viewed as OSs where there is also delayed commitment.)
- 11) *memory alloc & dealloc:* Memory deallocated by events is delayed until commitment in TW. No such feature of convention OSs.
- 12) *error handling:* In TW it is bound up with commitment. Not so in other OSs.
- 13) *I/O:* In TW it is bound up with commitment. Not so in other OSs.
- 14) *job termination:* In TW it is bound up with commitment. Not so in other OSs.